

**IN THE UNITED STATE PATENT AND TRADEMARK OFFICE
APPLICATION FOR UNITED STATES LETTERS PATENT**

**DISTRIBUTED EXTRACT, TRANSFER, AND LOAD (ETL)
COMPUTER METHOD**

By:

**Duane L. Porter
2813 NW Westbrook Circle
Blue Springs, Missouri 64015
Citizenship: United States**

**David E. Swanholm
16415 Nall Avenue
Stilwell, Kansas 66085
Citizenship: United States**

TITLE

Distributed Extract, Transfer, and Load (ETL) Computer Method

CROSS-REFERENCE TO RELATED APPLICATIONS

Not applicable.

**STATEMENT REGARDING FEDERALLY SPONSORED
RESEARCH OR DEVELOPMENT**

Not applicable.

REFERENCE TO A MICROFICHE APPENDIX

Not applicable.

FIELD OF THE INVENTION

[0001] The invention is a distributed extract, transfer, and load (ETL) computer method (hereinafter referred to as Distributed ETL) for linking together multiple information domains residing across an enterprise-wide computing environment in order to share corporate information, while reducing the time required to deliver that information. The invention overcomes the many-to-many relationships that exist between information sources and targets residing within each domain by implementing a central router to deliver information to and from these domains.

BACKGROUND OF THE INVENTION

[0002] Today, better integration of business information systems throughout the enterprise and incorporating more recent database management technologies results in real managerial competitive advantages. Information to support management decision making typically is extracted from operations systems and loaded into a data warehouse. Information as used herein includes data internal to a corporation that is generated in the day-to-day business operations of the corporation, aggregations and summaries of the internal data, and external information such

as demographic and geospatial data that can be related to the internal data. A data warehouse is a large storage facility within the information technology, IT, system that typically holds the current and historical information across an entire enterprise. It most likely consists of numerous databases and may exist on multiple platforms (operating systems). A database can be defined as a collection of information organized in such a way that a computer program can quickly select desired pieces of data. One can think of a database as an electronic filing system.

[0003] In the past, it was only viewed as necessary to capture periodic data extracts supplied by operational computing systems to update the data warehouse on a weekly or monthly basis. Strategic dependence on data warehousing was viewed as long-term and historical, not as a real time activity. In contrast, today's business environment is worldwide and does not stop at five o'clock P.M. For example, the Internet allows commerce to continue 24 hours a day, seven days a week. As a result, customer service support and normal telephony operations within many business organizations likewise operate around the clock.

[0004] Furthermore, marketing practices have evolved from mass media advertising to more personalized target marketing, which cuts expenses while providing customized treatment. To enable customized marketing, companies must know and understand their customers. Much of this understanding is made possible through groups such as Decision Support Services (DSS), the business group within an enterprise that has traditionally supported data warehousing. To be effective in supporting the extended operating and marketing efforts of an enterprise, information should be processed and pushed back from DSS into the customer contact areas such as customer service support, telemarketing and advertising campaign management as rapidly as possible. Thus, it is preferred that DSS be capable of receiving, processing, and providing information on an ongoing, near real time basis.

[0005] This invention, referred to as Distributed ETL, improves on current DSS services by providing the advantages of scalability of an extract, transfer and load (ETL) application and the near real time event processing of an enterprise application integration (EAI) application without the disadvantages associated with each.

SUMMARY OF THE INVENTION

[0006] The present invention is a distributed extract, transform and load (ETL) method for delivering information within a computing environment, comprising extracting information from an information source and transforming the extracted information. The transformed information is isolated by wrapping the transformed information into a message envelope having a standard format. The message envelope is routed to at least one information target, unwrapped to reveal the received information, preferably transformed again, and loaded into the information target. The extraction, transformation, and adaptation steps on the source side are isolated from the routing step such that the extraction, transformation, and adaptation steps on the source side may be executed simultaneously for a plurality of information sources distributed across the computing environment to produce a plurality of message envelopes. The routing, unwrapping, mapping, transformation, and loading steps on the target side are repeated for each of the plurality of message envelopes.

DESCRIPTION OF DRAWINGS

[0007] Figure 1 is a diagrammatic representation of business information domains.

[0008] Figure 2 is a diagrammatic representation of a basic hub and spoke architecture.

[0009] Figure 3 depicts the Distributed ETL architecture of the current invention.

[0010] Figure 4 is a diagrammatic representation of a staging area.

[0011] Figure 5 is a detailed diagrammatic representation of the router components.

[0012] Figure 6 depicts the utilization of the meta data repository by the adapter component of the Distributed ETL application and a parsing of the message envelope that is distributed by the router component.

[0013] Figure 7 is a more detailed diagram of the source and target spoke components.

[0014] Figure 8 is a diagrammatic representation of information flow with the Distributed ETL implemented on a single platform.

[0015] Figure 9 is a diagrammatic representation of a preferred embodiment of the Distributed ETL routing across multiple platforms.

[0016] Figure 10 depicts a general overview of how the Distributed ETL architecture integrates the various business enterprise information domains.

DETAILED DESCRIPTION OF THE INVENTION

[0017] Within the enterprise, information takes many forms depending on its usage. In other words, the business usage of the information drives its resultant structure. This information can be grouped into information domains according to its business usage or function as shown diagrammatically in Figure 1. The information domains include operations 10 used for daily operation information such as on-line transaction processing (OLTP); exploration zones 120 which includes data mining; data warehouse 40 which is the consolidated information store; and summarized views 150 used for predefined queries and on-line analytical processing (OLAP) reports.

[0018] As shown in Figure 1, the information domains overlap, indicating that the boundaries between business information domains are not absolutely crisp. Some warehouse information can be summarized, a characteristic of OLAP summary reporting, while OLAP tools sometimes drill down into detailed warehouse tables. Because of the differing data structures

required to support the differing business uses for the data, the same information may be used for multiple business functions. Therefore, the information domains are interdependent upon each other. The relationship between information domains can become a complex web linking information sources and information targets. These relationships can be intermixed and multidirectional, thus defining a many-to-many relationship between business information domains.

[0019] Current enterprise application integration (EAI) applications, such as IBM's MQ Series, Tibco's Rendezvous or Saga's Saga Vista, are available to coordinate information flow between information domains to facilitate business processes. These traditional EAI systems are synchronous, transaction driven systems that operate in real time or near-real time, but have certain other disadvantages (e.g., scalability limitations). An alternative tool is a standard extract, transform and load (ETL) application, which is an asynchronous, batch driven system that can be scaled up but typically suffers from delayed event processing.

[0020] Traditional EAI applications direct the flow of information between domains through a center point called a hub, or more commonly a message broker to manage many-to-many relationships in a scalable manner. The hub and spoke architecture is applicable to managing ETL needs that exist between information domains as well. Figure 2 is a diagrammatic representation of a basic hub and spoke architecture. The staging hub 70 is at the center of the hub and spoke figure with connecting spokes 2, 3, 4, and 5 to the various information domains, namely operations 10, exploration zones 120, data warehouse 40, and summarized views 150, respectively.

[0021] In contrast to the traditional EAI and standard ETL applications which rely on the basic hub and spoke architecture described above, this invention is a distributed extract, transfer,

5 and load (ETL) computer application to address the need for linking multiple information
6 sources and targets together to deliver shared corporate information, while reducing the time
7 required to deliver that information. The Distributed ETL overcomes the many-to-many
8 relationships that exist between information sources and targets by implementing a spartan
9 central router passing information to and from these domains through "thick spokes."

10 **[0022]** Figure 3 depicts the Distributed ETL architecture of the current invention.
11 Distributed ETL is comprised of two parts. First, a set of distributed "thick spokes" 20, 50, 100
12 and 130 which provide pre-routing information processing specific to sources connected by input
13 spokes and post-routing information processing specific to targets connected by output spokes.
14 Second, a spartan router component 75 serves as a central core by connecting the input and
15 output thick spokes. Preferably, all shared information in the enterprise will flow through the
16 router component.

17 **[0023]** The "thick" spokes 20, 50, 100 and 130 are used to manage specific source/target
18 information processing (also referred to as staging processes, which comprise extraction,
19 transformation, adaptation, and loading functions) as may be required prior to and/or after
20 routing of the information. In contrast, the router component 75 is spartan in that it primarily
21 manages delivery (i.e., routing) the information and does not perform processing steps thereon.
22 Staging processes that are specific to a source or target are pushed out closer to where the spoke
23 connects to that source or target, in contrast to being performed by the router component. By
24 isolating the router component to purely delivery functions, the entire architecture is streamlined,
25 enhancing its speed, efficiency, parallelism, and scalability. In sum, the architecture of the
26 present invention comprises: 1) the staging processes to manipulate the information (in the

“thick” spokes), and 2) the spartan router component to deliver the information. Therefore, the overall architecture can best be referred to as a Distributed ETL (extract, transform and load).

[0024] The staging processes, or the “thick” spokes of the Distributed ETL, are comprised of several software components to process information as necessary prior to and/or after routing it. Referring to Figure 4, a Distributed ETL architecture is shown wherein information is being routed from any one of two information sources 12 and 42 to any one of two information targets 122 and 152. Any operable number of information sources and targets can be linked using the Distributed ETL of this invention, and the number of sources and targets shown in Fig. 4 was chosen for convenience purposes only. On the information source side of the router component 75 (i.e., the source “thick” spoke), the Distributed ETL comprises extract components 22 and 52; transform components 26 and 56; and source adapter components 28 and 58. On the information target side of the router component 75 (i.e., the target “thick” spoke), the Distributed ETL comprises target adapter components 102 and 132; transform components 106 and 136; and load components 116 and 146. Additional adapter components 24 and 107 may optionally be used on the source spoke and the target spoke, respectively. Combined, these components insolate the information sources and targets from the router component and provide a standard interface for the router.

[0025] The first step, beginning at an information source locations 12 or 42 involves the extraction of information by the extract component 22 or 52. The extract component is a software program that functions as a connector to an external information source for retrieving information from the source. Information can be pulled or pushed from the source, and preferably the source pushes the information to the Distributed ETL environment. In most cases, it is desirable to retrieve only updates, revisions, changes and the like to the information

(“information changes”), as this limits the amount of information that needs to be delivered to the targets. Methods to extract the information changes from the source include:

1. A full extract of the source information is created, which is then compared against a previous full extract that was created at an earlier point in time. A detailed comparison between the two extracts generates the list of information changes.
2. A relational database management system (RDBMS) creates a log of all changes made to the database. This log can be read periodically to obtain the changes made to the source information.
3. The RDBMS can utilize what are called “triggers” that perform actions when certain data tables are affected and identified conditions occur, pushing the information changes to a table or queue that may be accessed for processing. The principal disadvantage to this approach is that it often loses transactional integrity. For example, where a single transaction for a customer spans several related tables, triggers looking for changes to only a single table may lose sight of the relationships needed to define all of the information changes.
4. The RDBMS can use a trigger to execute a “stored procedure” to capture information changes. In this case, the trigger executes when it senses a key relation change, and initiates a programmed stored procedure that gathers all information relevant to a logical transaction, based on the business rules for relationships contained in the referential integrity of the database design. The stored procedure then passes the information changes to the Distributed ETL with full transactional integrity.

[0026] The extracted information is then passed from extract component 22 or 52 to the transform component 26 or 56, respectively. Operational systems are the primary information

sources in an enterprise and typically employ highly normalized data structures that optimize OLTP performance critical to real time, customer facing processes. Normalized data structures tend to use multiple tables to minimize recurring data forms and data redundancy. Data redundancy and recurring formats within records/rows can cause severe performance penalties in an environment having a high number of information updates. The fragmentation of normalized design, useful in operational systems, creates problems when shipping the information in a normalized format to remote targets. The target must know the business rules, or referential integrity, that logically ties the several table structures together as an integrated, functioning whole. In order to pass the normalized data in a message format, each piece of information must be tagged in a way that relates it to the other associated information so they can all be delivered together and assembled according to the master blueprint to a state of referential integrity, which is very impractical. Thus, in order to implement message based transport, denormalization of the data to some extent is required so that sufficient information is contained in each message to allow interested targets to easily identify and subscribe to events they wish to receive. This is the essence of packet based messaging, i.e., that each packet is self-addressed by the information contained within it. The denormalization process required to create these packets relies on transformation steps, with full knowledge of the source referential integrity rules (i.e., business rules governing the information structure).

[0027] The transform component's function is to perform transformations on the extracted information. As described above, transformation refers to the application of specific business rules governing the structure of the information. For example, captured information such as the price of a certain product and the quantity sold to a customer may need to be transformed. The desired form of this information might be the extended price. By applying the business rule for

5 extended price, e.g., product price multiplied by quantity sold, the information has been
6 transformed to the desired structure. Another example of information transformation is where a
7 source captures customer names in the form of full name including the middle initial, e.g., John
8 D. Smith. But, the desired form of this information for one specific target could be first and
9 middle initials followed by the last name, e.g., J. D. Smith. Additionally, a different target may
10 desire the same information in yet another format, such as last name first, followed by first name,
11 followed by middle initial, e.g., Smith, John D. The last two examples demonstrate how two
12 different targets can access the same source information and transform it in different ways that
13 are acceptable to their business needs.

14 **[0028]** Transformations are used whenever the format of the information at the target, as
15 defined by the business needs, does not match the format of the information extracted or arriving
16 from the source. Stated alternatively, transformations are used to apply business rules that are
17 common to the enterprise on the source side, and specific to the target on the target side.
18 Typically, a source-side transformation is performed in order to facilitate packet based
19 messaging, as discussed previously. Also, a target-side transformation typically is performed in
20 order to facilitate the incorporation of information from multiple sources into a single target. In
21 other words, a benefit of information targets is that they consolidate information from multiple
22 sources, which frequently involves transformation. Even where a target receives information
23 only from a single source, typically the target adds value by creating different perspectives of
24 that information through aggregation and summarization, which also involve transformation.
25 Where the target only displays detailed data without further transformation, such information
26 may be read (i.e., queried) directly from an enterprise data warehouse, which typically has all
27 source rules already applied. However, a data warehouse typically uses a normalized data

5 structure, and thus a transformation may be required in SQL queries to return information in a
6 meaningful format.

7 **[0029]** The transformed information is now ready to be delivered to the target but must first
8 pass through the source adapter component 28 or 58. The source adapter component functions
9 asynchronously as a connector to the router component 75 thereby isolating transformations from
10 the router component and ensuring that the router component's primary function is delivery of
11 the information (in contrast to the pre-routing information processing in the "thick" source
12 spokes and the post-routing information processing in the "thick" target spokes). Herein lies
13 (both for the source adapter and as discussed below, for the target adapter) the adapter
14 component's primary function, to standardize and simplify the connection to the router
15 component. The source adapter component 28 or 58 is implemented in software and wraps the
16 transformed information in a standard message format, preferably a structured event format such
17 as CosNotification, which is part of the publicly available CORBA standard. The
18 CosNotification structured event is supported by various application server platforms such as
19 WebLogic available from BEA Systems, Inc. Hereinafter, the adapted information is referred to
20 as a message envelope.

21 **[0030]** The message envelope is sent to a message queue by the adapter to await distribution
22 via the router component 75. The source adapter must be intimately aware of content of the
23 extracted and transformed information given that a single source may provide multiple types of
24 information such as an addition, a revision, or a deletion of information. Because some targets
25 may wish to subscribe to, or receive, only some of these information types, the source adapter
26 must provide the specific information type for each message envelope sent to the router. To do
27 this, the source adapter must have access to the business rules that define how to identify each

type of information. For example, a certain code or date value may differentiate between an addition and a revision. Once the information type is identified, the source adapter tags the message envelope header with the proper message name (also referred to as the source ID) that will then be used for routing to the appropriate subscribing targets.

[0031] Meta data is used for control and management of information within the Distributed ETL application. Meta data is defined as “data about data” or facts about the information architecture within an organization. Meta data describes how, when, and by whom a particular set of information was collected, and how the information is formatted. Meta data is essential for understanding information stored in data warehouses, and the meta data is stored in a meta data repository. An example of meta data stored in the repository is transformation rules that describe how to move information from a transactional source database to a multidimensional target database and the aggregations and calculations that are applied to the information coincident to such a move. Another example of meta data in the repository is business rules (such as the definition of revenue or a common list of product codes) that multiple software components may want to access. Another example is when moving information from a source to a target, the meta data is checked to determine the appropriate characteristics of the data. For example, where a numeric field is required, an alpha numeric information piece would be rejected. When the adapter component detects such an error, the adapter may either route the information to an error queue rather than sending it on to the router component, or alternatively correct the information on the fly according to predetermined business rules. In sum, meta data provides transparent access to information.

[0032] Because the information processing within the thick spokes of the Distributed ETL runs in parallel, synchronization is important and can be accomplished by tracking the

information with meta data. The source adapter component utilizes a meta data repository to provide a “road map” and rules for the information that it processes. The meta data stored within the meta data repository can be referred to as the content definition. Analogous to a COBOL copybook in a COBOL program, the content definition specifies the structure and naming conventions for the content, that is how specific fields are names, field sizes, formats, etc. within the record content. The key associated with this content definition that defines where in the meta data repository the full content definition can be found is referred to as the content definition ID. These are the keys to the meta data and only the keys need to be passed through the router component rather than the whole content definition.

[0033] In a preferred embodiment, the target knows the source that it will be subscribing to in advance and simply references the content definition from the metadata repository when the target adapter is constructed. Thus, no content definition IDs (i.e., keys) need to be passed at all. Alternatively, where the content of each message may rely on a different document type definition (DTD) format, for example source information comprising extended markup language (XML) files, the key for the DTD may be sent with the message for retrieval on the target side rather than shipping the entire DTD with each message envelope. This significantly reduces the data transmitted and increases efficiency.

[0034] The utilization of the meta data repository by the adapter components of the Distributed ETL application and a parsing of the message envelope that is distributed by the router component are depicted in Figures 6. The source adapter component 28 pulls the content definition 61 from the meta data repository 60. The adapter then uses the content definition 61 to help process the transformed information (referred to as content 29) passed to it from the transform component as previously described in Figure 4. The content is the information

5 extracted from the source, expanded and/or modified by the transform, and its definition is stored
6 in the metadata store. Preferably, the content definition is defined and stored prior to the
7 processing of any information by the adapter. The adapter creates the message envelope (using a
8 structured event format) by performing simple mapping of the content to a uniform format and
9 by building header information such as the message name and source ID, which is used to route
10 the message envelope.

11 **[0035]** The adapter component attaches a header to the content creating a message envelope
12 30. The header includes a source ID 31 used to route the content to the appropriate targets by the
13 router component. The message envelope 30 may also contain the content definition ID 32 and
14 preferably will include the content 33. The router component delivers the message envelope 30
15 by utilizing a router table 86 as shown in Figure 5 and discussed below. The router table 86 is
16 stored in local memory to do an address lookup to find the subscribing target 35 and match its
17 source ID 34 with the source ID 31 in the message envelope 30.

18 **[0036]** Referring to Figure 4 on the exiting side of the routing function (i.e., the thick target
19 spoke), near the target or targets, additional adapters 102 and/or 132 implemented in software
20 read the router component output queue and remove the envelope from the distributed
21 information. Referring to Figure 6 the adapter 102 retrieves the content definition 62 from the
22 meta data repository 60 by reference to the content definition ID 32 of the retrieved message
23 envelope 30. Referring to Figure 4 transformations can be made in the target thick spoke (i.e.,
24 post-delivery by the router component), if necessary, according to the target requirements, with
25 an additional transform component 106 or 136. The information is then passed to the respective
26 load component 116 or 146. The load component preferably uses an Application Programming
27 Interface (API) software function to load the information into the target 122 and 152. An API is

5 useful to abstract the load operation from the specific characteristics of a given physical
6 datastore/file. In a typical embodiment, the target is a data warehouse and/or a data mart.

7 **[0037]** Multiple adapters are possible in the source and target thick spokes, as shown in
8 Figure 4. On the source side, additional source adapter 24 is shown between the extract
9 component 52 and transform component 56. On the target side, additional target adapter 107 is
10 shown between transform component 136 and load component 146. These additional adapters
11 allow for different technologies to coexist without requiring extensive and costly integration.
12 For example, an additional adapter may be used on the source side for transporting an extracted
13 file to a different platform for transformation. Another example, an additional adapter may be
14 used to enable a transactional system to flow events directly into a transformation process in
15 near-real-time (shown as process 7 in Figure 7).

16 **[0038]** The router component 75 is an important component in the Distributed ETL, because
17 all information sources and targets are connected through it. The key function of the router
18 component is to successfully route each message envelope from its publishing source to a single
19 subscribing target or multiple subscribing targets. It is comprised of multiple software
20 components that are message based, reading input from and writing output to message queues,
21 and should be generalized to the extreme. Ideally, routing should be easy to change and require
22 a minimal time to execute.

23 **[0039]** The router component 75 uses memory-based messaging to transfer information
24 between publishing sources and subscribing targets. More specifically, messaging occurs on the
25 source side between the source adapter and the router and on the target side between the router
26 and the target adapter. In summary, information is extracted from a source and transformed.
27 The adapter component creates the message envelope and publishes the source message to the

5 router via a message queuing mechanism. The router subscribes to all source queues, thereby
6 receiving the message envelopes to be routed. The router performs a lookup in the routing tables
7 using the source ID as a key. The router publishes the message envelope to one or more target
8 destination queues. The target adapter subscribes to a message queuing mechanism, unwraps the
9 message envelope, and passes the information to a target transformation component. The
10 information is then passed to a load component where it is written directly to the ultimate target.
11 Linking distributed sources or targets with messaging functionality enables a minimum handling
12 time per unit of information. Pairing transactional information extraction and loading with
13 messaging enables near real time processing of extraction and loading functions.

14 **[0040]** Guaranteed message delivery as a part of the publish and subscribe functionality is
15 important to ensure system reliability and verify that during exchange of messages no
16 information packets are lost. When the network is available, messages will be delivered
17 immediately. When the network is down, messages are still guaranteed to be delivered,
18 eventually, but with some latency. In the latter case, middleware guarantees that messages are
19 delivered as long as the network becomes available within a specified time period. In addition,
20 messaging should provide the assurance of non-duplicate delivery, meaning the messages are
21 only delivered once.

22 **[0041]** Messaging utilizes a queue system, which allows for a plurality of information
23 sources to simultaneously publish to a router component inbox queue, preferably where each of
24 the sources is capable of publishing independently of the others. A source (via the extract,
25 transform, and adapt components discussed previously) can publish to an inbox queue while the
26 router component is subscribing to that inbox queue. Unlike a sequential file, the router
27 component can read message envelopes from the front of the inbox queue at the same time that

the source is publishing new message envelopes into the back of the inbox queue. This allows faster aggregate throughput than a traditional batch environment. A router component can then publish a message envelope to a subscribing target's queue, which the target can access (as discussed previously with regard to the target adapter) and store for later retrieval matching its business cycle needs.

[0042] Performing its functionality in memory keeps the router component robust, as memory is hundreds to thousands of times faster to access than disk. Thus, it is preferred to use routing tables that can be maintained in memory for the control of routing destinations. The routing tables are backed up by persistent storage such that the tables may be reloaded to memory when the power is restored after a failure, or to update the content of the table when target subscriptions are changed.

[0043] Figure 5 is a diagrammatic representation of the router components and functionality comprising the Distributed ETL architecture. The router component 75 receives a message envelope 82 from any one of multiple sources 12, 13, 14, 182, 183 or 184 through the thick source spokes and the inbox queue discussed previously. The message envelope is then passed to an address lookup 84 where a publish/subscribe cross-reference table 86 (also referred to as a router table) is accessed to determine where the message envelope is to be sent. The header on the message envelope identifying the publisher (i.e., source) corresponds to a list of subscribers (i.e., targets) in the cross-reference table 86 so the router component can determine the routing destination. The message envelope is then sent from Address Lookup 84 to Envelope Re-Address 94 where the target address is added to the message envelope. In an alternative embodiment, the source adapter automatically sends the message envelope to the router, which in turn routes the message envelope based upon the originating source (via the lookup table) and

with the original source name intact. In other words, the target address is found in the lookup table and is not placed upon the message itself. The message envelope is then sent out to the appropriate target via Envelope Transmit 96. The output from the router component is written to an output queue that will be read by an adapter component on the target side of the distribution function. Persistent Publish/Subscribe Tables 90 are updated through a Publish/Subscribe Manager 92, which provides a Graphic User Interface (GUI) configured for table management. These updates are sent to the Memory Refresh 88 to provide current entries in Publish/Subscribe cross-reference table 86. Addresses contained within the cross-reference table are stored in local memory within the router to quicken address lookup time. The Error Handle and Administration 97 handles errors by sending an error message, Alert Message Handling 98 and storing the errors in a Routing Error Cache/Queue 99.

[0044] Downstream of the routing function another adapter 102 reads the message envelope 30. It uses the content definition 62 retrieved from the meta data repository 60 to extract the content 101 from the distributed information, as described previously. The information is subsequently processed in the thick target spoke on its way to the appropriate target as described previously in relation to Figure 4.

[0045] Figure 7 is a more detailed diagram of the components comprising the thick source and target spokes. As described previously, information is extracted from the information source 12 or process 7 by the extract component 22. The extracted information is passed to the transform component 26 where the information is transformed. The extracted, transformed information is sent through the adapter component 28 where it is wrapped in a message envelope and sent to the inbox queue for delivery by the router component 75.

[0046] On the output side of the router component 75, the adapter component 102 reads the information from a message queue and takes the information out of the message envelope (removes the header). In a preferred embodiment, the source information (including the type of information such as an addition, deletion, etc.) contained within the message envelope is specific enough such that the router component can route the message envelope only to targets desiring (i.e., subscribing) such information, thereby eliminating the need for filtering. In a non-preferred embodiment, the adapter component 102 passes the information through an optional filter component 104. The filter component's function is to only accept information required by the target and to reject non-conforming information. The filtered information may now undergo more transformations by the transform component 106 that can be applied at the event level, i.e., one message at a time. In other words, all information need not be made available at the same time such as sorting a group of records, all of which must be made available in order to perform the search. In contrast for example, a trunk code could be used in a message to perform a lookup and the trunk location returned to update the same record. Where the information has been acquired continuously over a period of time, the information in a preferred embodiment is passed to the time-period aggregator component 108. Here the information can be accumulated in an "active" partition or file and later retrieved by toggling the "active" partition 110 to an "inactive" state 112, allowing it to be read into the target transformation. The information can be held to accommodate a business time cycle, e.g., weekly or monthly reporting. When the partition is toggled to an inactivate state 112, the accumulated information can now be passed on for additional transformations by an additional target transform component 114 where necessary. This level of aggregation (in contrast to a time period aggregation) accommodates business aggregations, e.g. by region, product, salesperson, etc. as well as other transformation logic such

as any matching operations if multiple information sources are received for consolidation. The information is then passed to the loading component 116 where the information is finally loaded into the target 122 typically as a batch process.

[0047] Figure 8 is a diagrammatic representation of information flow on a single platform of a Distributed ETL. A single platform means that all the information sources and targets are routed through one router component. The single platform could represent a small company's information system or a portion of a larger company's information system. There can be multiple sources 12, 13 and 14 (connected to the router component by thick source spoke 330) and multiple targets 122, 123 and 124 (connected to the router component by thick target spoke 335). Note that targets can also be sources 182, 183 and 184 known as recursive targets, both receiving information through a thick target spoke 340 and sending information through a thick source spoke 340. A recursive target receives information from the router, adds some unique transformation value to the information and sends it back through the router to another specific target or targets. A recursive spoke does not send out the same information it receives. Information flow begins at the information source 12, 13 or 14 travels through thick source spoke 330 to the router component 75 and then on through thick target spoke 335 or 340 to either the information target 122, 123 or 124 or to the recursive target 182, 183 or 184, respectively. The uniquely transformed values from the recursive targets are then routed back through thick source spoke 340 to the router component 75 and on through thick target spoke 335 to the information target 122, 123 or 124.

[0048] In contrast to the single platform structure of Figure 8, Figure 9 is a diagrammatic representation of a preferred embodiment of the Distributed ETL routing environment across multiple platforms. A plurality of sources and targets 12-14, 42-47, 182-184, 122-124, 152-157,

282-287 within an enterprise can be logically grouped into multiple platforms as shown with platform A 11, platform B 41, and platform C 121. By introducing multiple routers 75, 76 and 77 into the overall architecture which deliver message envelopes to one another a represented by arrows 300, 305, and 310, better utilization of bandwidth between devices can further enhance efficiency as well as make the design easier to scale up to very large business intelligence integration systems.

[0049] Integration of the various business enterprise information domains from Figure 3 with the Distributed ETL architecture of Figure 9 is depicted in Figure 10. The information domains 10, 40, 120 and 150 link to a logical hub 70 to resolve their many-to-many relationships. Platforms A 11, B 41, and C 121 show the physical means by which this is accomplished within the logical hub. The physical platforms contain multiple source and target spokes, which could connect to any of the information domains. Information is shared between physical platforms through the routers that exist on each platform 75, 76, 77. Figure 10 shows how the physical solution fits into the larger logical/conceptual picture.

[0050] In a typical embodiment, the information target comprises a data warehouse and a data mart. Since it can be quite cumbersome to access information from the data warehouse because of the large amounts of data and typically normalized data structures, a common IT design feature is to provide smaller information storage places that focus on a particular subject or department, called data marts. Historically, these can either be dependent data marts, which are subsets of the data warehouse, or independent data marts, which load information directly from operations (point to point). Data marts function to provide easier access to stored information. The data marts are highly indexed and optimized for query. Where one knows in advance the questions to be asked, data marts can be structured accordingly.

5 [0051] Legacy systems utilizing dependent data marts tend to be slow because it takes time
6 to first load the data warehouse and then initiate another process to load the data marts. To
7 reduce latency, an efficiency feature of the Distributed ETL invention loads information directly
8 into 'co-dependent' data marts instead of the data warehouse. Co-dependent data marts are built
9 using information that is dependent upon the same rules as the data warehouse but does not
10 require the latency involved in writing to a data warehouse structure followed by extraction,
11 transport and/or transformation and a final write to the dependent data mart.

4000-04200-1616